# Paying for QoS: An Optimal Distributed Algorithm for Pricing Network Resources

Errin W. Fulp[*†], Maximilian Ott[*], Daniel Reininger[*] and Douglas S. Reeves[†]

## Abstract

Network applications require certain individual performance guarantees that can be provided if enough network resources are available. Consequently, contention for the limited network resources may occur. For this reason, networks use flow control to manage network resources fairly and efficiently. This paper presents a distributed microeconomic flow control technique, that models the network as competitive markets. In these markets switches price their link bandwidth based on supply and demand, and users purchase bandwidth so as to maximize their individual Quality of Service (QoS). This yields a decentralized flow control method that provides a Pareto optimal bandwidth distribution and high utilization (over 90% in simulation results). Discussions about stability and Pareto optimal distribution are given as well as simulation results using actual MPEG-compressed video traffic.

## 1 Introduction

Advances in computer network technology have resulted in complex networks that must accommodate a variety of network applications. These applications transmit a range of information, from simple text and graphics to complex interactive voice and video. Each application requires a certain Quality of Service (QoS), which may include bounds on the packet delay, delay variation and loss rate. These service guarantees can be provided if the network resources are available, such as link bandwidth, buffer space and processor time. Since the amount of resources is finite, contention may occur. For this reason, networks need a method of flow control to manage resources in a fair and efficient manner.

There are two goals associated with flow control, fairness among applications and the balance between throughput and QoS [2, 7]. Defining fairness is difficult because of the various types of applications and their desired QoS. The balance between throughput and QoS is the concept that the network should seek high resource utilization, but not at the expense of poor QoS (and vice versa). Hence, due to heterogeneous networks, diverse resource requirements and the goals associated with flow control, proper flow control is a challenging problem. Several different methods of flow control have been proposed, some specifically for certain types of networks. We will briefly discuss the general classes of flow control as well as a new type based on economics.

Preventive flow control determines the transmission rate of each source that will avoid congestion. In this case congestion is prevented and some service guarantees can be provided. However, this type of flow control may lead to over allocation of resources and is not well suited for the dynamic changes (such as variable bit rate sources) that may occur in the network. Feedback flow control methods alter data transmission to adapt to changing network conditions. Window flow control is one example used in packet networks. In this strat-

---

[*]C&C Research Laboratories, NEC USA, `max|djr@ccrl.nj.nec.com`

[†]Departments of ECE and CSC, North Carolina State University, `ewfulp@eos.ncsu.edu`

egy, network feedback is used to limit the number of packets transmitted; however this type of flow control is not well suited for large networks because of propagation delays and few (if any) QoS guarantees can be made [2]. In ATM networks, several feedback traffic management strategies have been proposed for Available Bit Rate (ABR) service. These traffic management techniques use network feedback to alter the rate of a source (instead of the number of packets). Examples of explicit rate techniques include EPRCA and ERICA [1]. These strategies rely on the circulation of a Resource Management (RM) cell per connection [1]. As the RM-cell travels along the path, a switch and/or the destination may alter its contents. Exactly how this is done depends on the strategy. Once the cell reaches the destination it is returned to the source, which must alter transmission based on the RM-cell information. When a switch becomes congested, these traffic management strategies seek to allocate the bandwidth in a *fair* (max-min) manner. However, these methods do not take into account the fact that some sources may be able to reduce their transmission rate (for example compressed video) more easily than others. Therefore when congestion occurs, this socialistic allocation may not be the best when considering the individual QoS expected by each user.

An economic flow control method models the network as an economy, then applies microeconomic principles for resource allocation. A simple network economy consists of two types of agents: consumers (network applications) and producers (switches). Consumers require resources to satisfy their QoS. Producers own the resources sought by consumers, and seek to maximize their satisfaction by selling or renting their resources. Using this framework, microeconomics can be used to define how network resources are allocated.

One approach of applying microeconomics to computer networks involves a maximization of utility functions [9, 10, 11, 12, 14, 17]. A utility function maps a resource amount to a satisfaction value. Using this function, one can compare the satisfaction levels of different resource amounts. The maximization process determines the optimal resource allocation such that the utility of a group of users is maximized subject to budget and resource availability constraints. Accurately grouping users together may be problematic due to the wide variety of applications and their diverse resource requirements. Another problem is that these approaches generally require a centralized entity to determine the optimal allocation amount. This is undesirable because the economy relies on one entity, which is not reliable or fault tolerant.

Another microeconomic approach, congestion pricing, charges users for their consumption of resources and resources are priced to reflect supply and demand [3, 4, 13]. With such a model, prices can be set to encourage high utilization of network resources as well as a fair distribution. Users act independently, attempting to maximize their own utility and prices are set based on local resource conditions. It has been shown that pricing based on supply and demand results in higher utilization than traditional flat (single) pricing [3, 13]. Ferguson, et al. is an example of flow control based on pricing network resources [4]. Prices of links in the system were iteratively adjusted until an equilibrium of supply and demand was reached. They were able to prove that the system achieved a Nash equilibrium; yet they required demands to be constant until the equilibrium price was determined. If the demands changed, the prices were no longer valid. Our approach uses congestion pricing in a competitive market. Similar to other microeconomic flow control methods, our approach is decentralized, seeks an equilibrium price and achieves a Pareto optimal distribution. In addition, our approach maximizes individual QoS, adapts to network dynamics and is scalable to heterogeneous networks.

The remainder of this paper is structured as follows. Section 2 reviews the competitive market model. Section 3 describes the pricing technique in detail. Section 4 discusses how our

pricing strategy achieves an equilibrium price and a fair Pareto optimal distribution. Section 5 discusses how the pricing policy contends with network dynamics such as, users entering/exiting and multimedia traffic. Section 6 describes the simulation results. Finally, section 7 reviews the pricing technique, summarizes the results and discusses some open questions.

# 2    Competitive Market Model

We will use a competitive market model for our network economy. The competitive market model consists of scarce resources and two types of agents, consumers and producers. A resource is an item (or service) which is valued by agents in the economy. Since it is scarce, there is never enough of the resource to satisfy all the agents all the time. For this reason, allocation decisions must be made. Consumers require resources to satisfy wants. Producers create or own the resources sought by consumers. These agents come together at a market, where they buy or sell resources. Usually these exchanges are intermediated with money and the exchange rate of a resource is called its price. Prices are set with respect to supply and demand. The price increases if the demand is greater than the supply and decreases when the demand is less than the supply. When they are equal, the market and price is in equilibrium. This moment is referred to as "clearing the market" and the resulting allocation is Pareto optimal [18]. Pareto optimality is the allocation of finite resources such that no sub-set of users can improve on their allocation without lowering the utility of another. This model was chosen for our computer network economy because of its ability to achieve certain desirable goals, such as Pareto optimal distribution and price stability. The competitive market also has a simple structure and a well founded mathematical basis for analysis.

# 3    A Proposed Pricing Policy

This proposed flow control method is based on a competitive market model, where pricing is done to promote high utilization and Pareto optimal distribution. There are three entities in this network economy: users (those who execute network applications), Network Brokers (NB) and switches. Using the competitive market nomenclature, users are consumers, switches are producers and network brokers are used to assist the exchange of resources in the market. While there are many resources in a computer network, this paper focuses on the pricing of link bandwidth.

## 3.1    Switch

In our competitive market, the switch owns the link bandwidth that is sought by consumers. The network consists of several switches interconnected with links. For a unidirectional link between two switches, we consider the sending switch as owner of the bandwidth of that link. Each switch prices its link bandwidth based on local supply and demand for that link. Therefore a single switch, having multiple output links, will have one price associated with each output port. For example in figure 1, switch 0 owns and will determine a price for link 0. The entire network can be viewed as multiple competitive markets, one market per link (similar to the New York Stock Exchange). These markets operate independently and asynchronously since there is no need for market communication (for example, price comparisons) or synchronization from switch to switch. Consequently, this results in a decentralized economy, where the physical failure of one switch/link does not necessarily cause failure of the entire economy.

The price computation for link $i$ is performed at the switch, at discrete intervals of time. We denote the $n$th calculation instant as $t_n^i$ and the interval of time between the calculation points $t_n^i$ and $t_{n+1}^i$ as the $n$th price interval, $P_n^i$. The price during $P_n^i$ is constant and is denoted as

$p_n^i$. The demand for bandwidth at link $i$ is measured as the total (aggregate) traffic received at its associated output port. During the $n$th price interval, $P_n^i$, the total demand is expected to change; even so, the calculation of $p_{n+1}^i$ will only use the demand measured at the end of the interval. For this reason, let the demand for bandwidth at link $i$, at the end of the $n$th price interval, be denoted as $d_n^i$. The supply of bandwidth at link $i$ is constant and denoted as $S^i$.

At the end of the price interval, $P_n^i$, the switch updates the price of link $i$ using the following equation,

$$p_{n+1}^i = p_n^i + c \cdot \left( \frac{d_n^i - \alpha \cdot S^i}{\alpha \cdot S^i} \right) \qquad (1)$$

The form of the price equation is referred to as a tâtonnement process and is used in a competitive market to set the price with respect to the current supply and demand [19]. In a tâtonnement process the new price is equal to the previous price plus a correction function. The correction function provides feedback based on the demand (received traffic) and the supply (bandwidth available). The bandwidth available is the total bandwidth times a constant $\alpha$, where $0 < \alpha \leq 1$. This causes the price to increase after some percentage ($\alpha$) of the total bandwidth has been reached. This is evident from the equation, since the price will only increase if the numerator is positive ($d_n^i > \alpha \cdot S^i$). The price will decrease as the demand decreases and will increase as the demand increases. An *equilibrium price* $p_*^i$ is reached at link $i$ when the supply equals the demand. At this point the market clears for link $i$ and the allocation of bandwidth is Pareto optimal [18]. The positive constant $c$ amplifies the feedback signal and its value ultimately controls how quickly the price will increase or decrease (speed of adjustment). Note that the equation can yield negative prices. We will assume that the price will not fall below a certain non-negative minimum price (set by the switch).

After the new price, $p_{n+1}^i$, is calculated, a new price quote is forwarded to each NB using this link. The price quote for link $i$, denoted as $q_{n+1}^i$, consists of; $p_{n+1}^i$, $d_n^i$, $S^i$, $c$ and $\alpha$. The NB will use all of the information in the price quote to determine the amount of bandwidth to purchase.

## 3.2 User

The user, executing a network application, requires bandwidth for transmission. The amount of bandwidth desired is determined from the application and is denoted as $b_m$. We assume $b_m$ is constant for the duration of the application. In section 5 we will allow $b_m$ to vary over time, which is desirable for multimedia transmission.

Based on prices and wealth, the user can afford a range of bandwidth (less than or equal to $b_m$), and some amounts will be preferred over others. In economics these preferences are represented with a utility function. The utility function maps a resource amount to a real number, that corresponds to a satisfaction level. Assuming $U(\cdot)$ is a utility function, if the user prefers an amount $x$ over $y$ (this is represented using the notation $x \succ y$) then $U(x) > U(y)$. The utility curve can be used to compare resource amounts based on the satisfaction the user will receive. This provides an important link between resource amounts and user satisfaction. For this economy we will use *QoS profiles* for the utility curves. Based on psycho-visual experiments, the QoS profile is a two dimensional graph, as seen in figure 2. The profile can be approximated by a piece-wise linear curve with three different slopes. The slope of each linear segment represents the rate at which the performance of the application degrades when the network allocates a percentage of the desired bandwidth ($b_m$). The horizontal axis measures the bandwidth ratio of allocated bandwidth to desired bandwidth ($b_m$). The vertical axis measures the satisfaction and is referred to as a QoS score. Our QoS scores range from one to five, with five representing an excellent perceived quality and one repre-

senting very poor quality. As seen in the figure, if the allocated bandwidth is equal to the desired bandwidth ($b_m$), the ratio is one and the corresponding QoS score is 5 (excellent quality). As this ratio becomes smaller the QoS score reduces as well. Profiles can be created for a variety of applications and redefined as users gain more experience. New and updated profiles can be easily incorporated within the economy as they become available. More information about QoS profiles is given in [15].

Finally, the user is charged continuously for the duration of the session (analogous to a meter). To pay for the expenses, we will assume the user provides an equal amount of money over regular periods of time. We will refer to this as the budget rate of the user, $W$ ($/sec). A single initial endowment could have been used, but would necessitate defining how it is spent during the session. To simplify simulation and analysis, budget rates are used.

## 3.3 Network Broker

Users can only enter the network economy through a network broker (NB). This entity is an agent for the user and is located between the user and the edge of the network. Representing the user in the economy the NB performs the following tasks: connection admission control, policing, and purchase decisions. Although the NB works as an agent for the user (making purchasing decisions), we assume that the NB operates honestly in regards to both the switches and the user.

The NB controls network admission by initially requiring the user to have enough wealth to afford at least an *acceptable* QoS; otherwise, the user is denied access. The purpose of this requirement is to be certain all users are viable consumers in the market and to prevent overloading the economy. We believe the social welfare of the economy is better when it consists of fewer users each receiving a good QoS, instead of many users each receiving a poor QoS. Hence, we are attempting to maximize the number of users in the economy, where

each user can afford an acceptable QoS. If the desired bandwidth is constant, then the test is relatively simple. However, for sources where the desired bandwidth will change over time, a more complex admission test is required.

The NB monitors the user and the prices by gathering and storing information about each. From the user, the NB collects and stores; the QoS profile, $b_m$ and $W$. The NB also stores the route, $R$, that connects source to the destination, where $R$ consists of $v$ links, $\{l^i, i = 1 \ldots v\}$. For each link on $R$, a price quote, $q^i$, is collected, where $\vec{q} = \{q^i, i = 1 \ldots v\}$ is the vector of price quotes for the route. Price quotes will change over time, since they represent link supply and demand. The NB will only store the most recent price quote from each link in the route. The NB will divide the budget rate, $W$, into a vector of $v$ budget rates $\vec{w}$, where $\vec{w} = \{w^i, i = 1 \ldots v\}$ and $w^i$ corresponds to link $i$. Separate budgets are used to localize the effect of prices to each link. This prevents spending the entire budget on one expensive link. Of course depositing and withdrawing to and from these individual budgets is possible and perhaps advantageous. This is one area for future work and it is not considered here. Using this information the NB levies the user for their consumption. Users will be charged based on usage (similar to electricity), since bandwidth is a non-storable item. Using this information the NB polices the user, ensuring only the bandwidth purchased is used.

Finally, the NB determines the amount of bandwidth to purchase. This value is based on the budget, current prices and QoS profile of the user. Denote the $r$th amount of bandwidth to purchase (use) as, $u_r$. Once the NB determines $u_r$, the user will start sending at this rate immediately. There is no need for direct confirmation/feedback from the switches. A new amount of bandwidth to purchase, $u_{r+1}$, will be determined in response to a new price (or change in demand, as will be described in section 5). Exactly how the NB determines $u_{r+1}$ is described next.

### 3.3.1 Determining the Bandwidth to Use

When determining $u_{r+1}$, the NB will first calculate the maximum and minimum bandwidth that can be used. The maximum bandwidth that can be used at link $i$ is,

$$b_{max}^i = \frac{w^i}{p^i} \,, i = 1 \ldots v$$

therefore the maximum bandwidth the user can afford is,

$$\widehat{b}_{max} = \min_{i=1\ldots v} \{b_{max}^i\} \,.$$

Note this equation maximizes the bandwidth at the current prices. The minimum bandwidth that can be used is determined from the QoS profile, $b_m$ and the value that corresponds to the lowest acceptable QoS score. It is possible that $\widehat{b}_{max} < b_{min}$ (the minimum is not affordable), due to the QoS constraint, prices and budgets. If this case arises, the user must either; increase the budget rate, accept a lower QoS, or drop the connection. Properly managing such a situation is an area for future work.

After $\widehat{b}_{max}$ and $b_{min}$ have been calculated, $u_{r+1}$ can be determined. The following procedure will attempt to find the maximum bandwidth at the current prices and budgets. It also calculates the price impact of the change in consumption on itself. In microeconomics this is similar to *internalizing externality*. The initial $u_{r+1}$ is,

$$u_{r+1} = \begin{cases} \min\{\widehat{b}_{max}, b_m\} & \text{if } \widehat{b}_{max} \geq b_{min} \\ \emptyset & \text{otherwise} \end{cases}$$

$$(2)$$

Using the price quotes, the NB must determine if $u_{r+1}$ will cause a price change that the user cannot afford, minimizing the externality of the bandwidth used. The highest price that the user can afford at link $i$ is,

$$\frac{w^i}{u_{r+1}} \,. \qquad (3)$$

The new price caused by $u_{r+1}$ at link $i$ is,

$$p^i + c \cdot \left( \frac{u_{r+1} - u_r + d^i - \alpha \cdot S^i}{\alpha \cdot S^i} \right) \,. \qquad (4)$$

The new price given in equation 4 can not exceed the maximum price affordable, given in equation 3. Using these equations the following inequality provides a bound on feasible $u$ values,

$$w^i \geq u_{r+1} \cdot \left[ p^i + c \cdot \left( \frac{u_{r+1} - u_r + d^i - \alpha \cdot S^i}{\alpha \cdot S^i} \right) \right]$$

$$(5)$$

Solving (5) for $u_{r+1}$ yields the bandwidth at link $i$ whose price change the user can afford. The inequality (5) has a closed form or it can be solved iteratively.

As described earlier, once the NB has determined its $u_{r+1}$ it will start sending immediately at this rate. No signaling is performed. This technique provides a significant reduction in overhead; however an over allocation of resources may occur. Consider the following scenario. Assume many users are using one link and the price has reached an equilibrium value. Now assume one user ends their session and this reduction of bandwidth results in a lower price. If the remaining users react to this lower price, over-allocation of bandwidth may occur. One simple approach to prevent this situation is to have the switch adjust $c$ so the price decreases at a slower rate. An over-allocation may still occur if many users using a link start sending at a higher rate simultaneously due to their application (not price); however this would require a correlation of these events. In general, adjusting the price based on $\alpha \cdot S^i$ and the high capacity of most links diminish the significance of this problem.

## 4 Optimality

As with any allocation strategy there are certain optimal allocation goals. Since pricing is used, optimality will be described in microeconomics terms. There are two important goals

this technique strives for; Pareto optimal allocation and price stability.

As described in section 2, Pareto optimality is the allocation of finite resources such that no sub-set of users can improve on their allocation without lowering the utility of another, given that supply equals demand. This is a standard goal in microeconomics for social benefit of resource distribution. Several proofs have been developed to show that competitive markets reach a Pareto optimal distribution [18]. A proof that our computer network economy achieves a Pareto optimal distribution is given in [5].

The equilibrium price ($p_*$) occurs when a price is reached such that the demand equals the supply. At this point, the resources are fully utilized. If the demand changes, pricing mechanism should alter the price to return to equilibrium. This property is what is referred to as price stability. A proof that our proposed pricing technique has price stability is also given in [5].

# 5  Network Dynamics

Thus far, the description and analysis of the network economy has not considered the dynamic nature of an actual computer network. The dynamics we are interested in include; users entering/exiting the network, and allowing Variable Bit Rate (VBR) sources. Although prevalent in actual networks, these dynamics have been either or both excluded in other microeconomic flow control methods.

As described in the introduction, multimedia applications will constitute a large portion of the applications in current computer networks. The traffic generated by these applications can be described as VBR, which means the bandwidth required will change often and unexpectedly. Restricting the user to a constant desired bandwidth, as described in section 3.2, requires the user to purchase the highest amount of bandwidth expected (peak rate). For VBR sources, this approach is both difficult to implement and inefficient. Implementation is dif-

ficult since the peak rate may not be known in advance (consider live or interactive video). Purchasing only the peak rate is inefficient since the application may only require the peak rate for a short period of time. For these reasons it is advantageous to allow the user to change the desired bandwidth over time. For a particular application, denote the $m$th desired bandwidth change as $t_m$, and the interval of time between bandwidth changes $t_m$ and $t_{m+1}$ as the $m$th application interval, $A_m$. The bandwidth desired during $A_m$ is constant and is denoted as $b_m$. It is important to note the length of $A_m$ depends on the application and will vary over time. At the end of $A_m$ the new desired bandwidth $b_{m+1}$ is sent to the NB. Now the NB determines a new amount of bandwidth to use, $u_{r+1}$, when either a new price or new desired bandwidth is received. The procedure for determining $u_{r+1}$ is described in section 3.3.1. Once $u_{r+1}$ has been determined the user starts sending at this rate immediately.

Since the number of users and demands for bandwidth change over time, the aggregate demand, $d_n$, for a link will vary as well. This can be depicted by shifting the demand curve (for the bandwidth of a link) left or right over time. As a result there is not a single equilibrium price, $p_*$, for all time. However, the market can be viewed as having multiple equilibrium prices, each for some segment of time. During a segment the pricing technique will seek the equilibrium price as described in section 4. Once this price is found, the resulting distribution is Pareto optimal. When the aggregate demand changes, the stability of the price equation ensures that the price of bandwidth always moves towards $p_*$.

# 6  Experimental Results

In this section the performance of the network economy is investigated via simulation. Previous microeconomic flow control techniques either do not provide experimental results or simulate limited networks (network size and/or traffic source types). Experiments performed
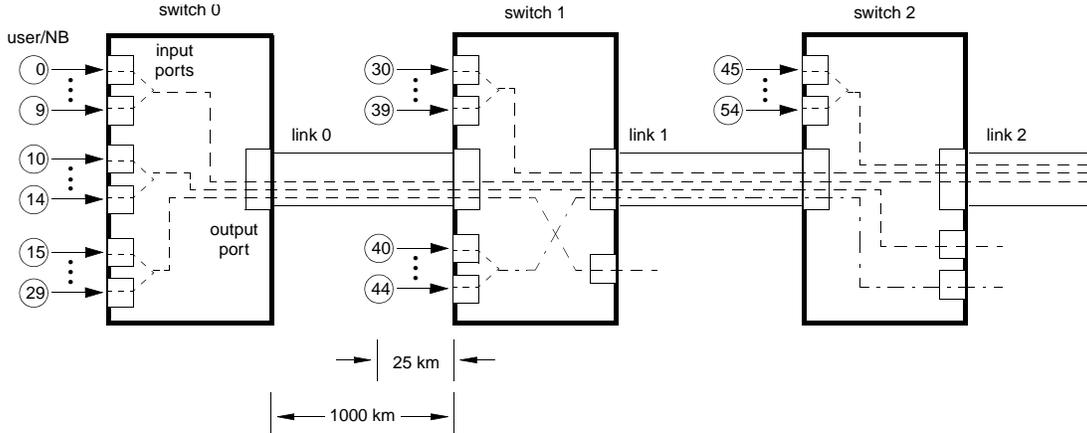
**Figure 1:** Network configuration used in simulations.

will consist of a realistic network configuration, allow users to randomly enter the network and use actual MPEG-compressed traffic. Experimental results will show that the proposed pricing technique achieves a fair Pareto distribution as well as high network utilization.

The network simulated consisted of 55 users/NB, three switches and three primary links, as seen in figure 1. Each output port carried traffic from 30 users and connected to a 55 Mbps link. Links interconnecting switches were 1000 km in length, while links connecting sources to their first switch were 25 km in length. Users had one, two or three hop routes and entered the network at random times, uniformly distributed between 0 and 60 seconds. The network can be described as a "parking lot" configuration, where multiple sources use one primary path. This configuration was agreed upon by members of the ATM Forum for allocation comparisons since it provides competition among users with different routes and various propagation delays [8].

The pricing strategy had the following initial values. Each user had a budget rate, $W$, of $3 \times 10^7$/sec [1]. Since all users have the same budget rate, they are considered equal (purchasing power). This should cause all users to be treated fairly, with no disproportionate allocation if all require the same amount. Switches initialized their prices to 1. The price equation $c$ constant was set to 50 and $\alpha$ (the target utilization) was 90%. We also assumed no propagation delay between the user and their NB. Switches updated their link prices at an interval ($P_n$) equal to 20 times the longest propagation delay of any user connected to it.

Each user (source) used the QoS profile given in figure 2, which was generated from an actual MPEG video application [15]. The source for each user was one of 15 MPEG-compressed traces obtained from Oliver Rose at the Univer-

---

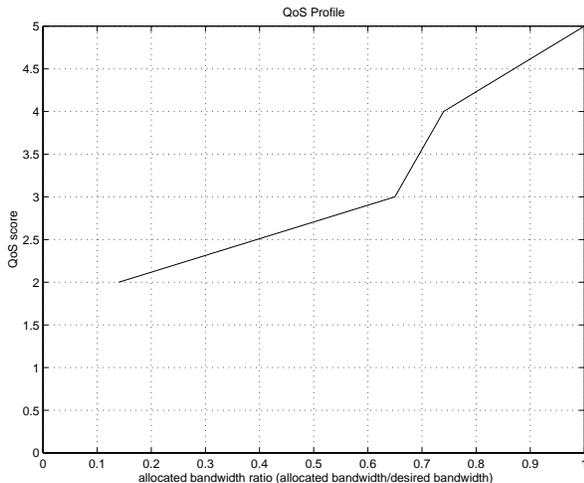[1]The denomination is based on bps, if based on Mbps the budget would be 300/sec.

**Figure 2:** QoS profile used in simulations.

sity of Würzburg, Germany [16] [2]. Each trace is a thirty minute segment of the original video and each was encoded with constant quality using the same MPEG-1 encoder card. Relevant statistics of each video are presented in [6] and [16]. As reported in [16], the Hurst parameters indicate all videos exhibit long-range dependency, and significant peak-to-mean ratios ranging from 18.4 to 4.63 based on average frames; therefore it is evident that these are very difficult sources to regulate. To date no other microeconomic flow control method has provided experimental results with actual MPEG sources.

We are interested in the link bandwidth utilization, the QoS provided to each user and the allocation optimality. Allocation graphs are provided to measure the utilization of link bandwidth. To quantify the QoS observed the percentage Good or Better (GoB) was calculated. This measurement is the average percentage of time a user had a quality score of at least 3. Finally, the optimality of the al-

---

[2] Traces can be obtained from the ftp site `ftp-info3.informatik.uni-wuerzburg.de` in the directory `/pub/MPEG`

location is given in the *fairness index* graph, which indicates how far the allocation is from optimal [8]. Suppose the allocation among $n$ users is $\{x^1, x^2, \dots, x^n\}$ and the optimal allocation (Pareto) is $\{\hat{x}^1, \hat{x}^2, \dots, \hat{x}^n\}$. Define the normalized allocation as $\tilde{x}^i = x^i/\hat{x}^i$ for each source, then the fairness index is computed as,
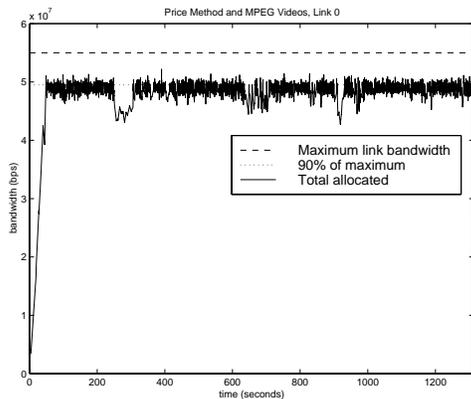
$$\text{fairness index} = \frac{(\sum \tilde{x}^i)^2}{n \sum (\tilde{x}^i)^2}$$

and will be plotted as a function of time. A fairness index of 1.0 indicates an optimal allocation while 0 indicates an unfair distribution. A measurement equal to or greater than 0.99 will be considered optimal [8].
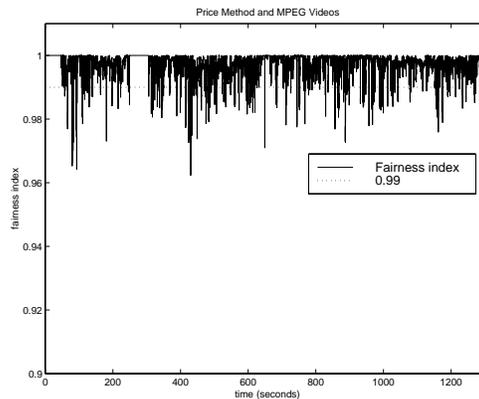
For this simulation, the bandwidth allocation for link 0 (representative of the other links) is given in figure 3(a) and the fairness index (entire network) graph is given in figure 3(b). The allocation graph shows that the total allocation of bandwidth stayed in the vicinity of 90% ($\alpha$, the target utilization). The fluctuation around this value is the result of users entering the network and changing demands. The GoB (for the entire network) was 98.6%, which means 98.6% of the time users measured a QoS score of 3 or better. The optimality of the allocation can be seen in the fairness index graph, where 92.7% percent of time the fairness index was 0.99 or greater, indicating an optimal allocation. During the remaining 7.3% of the time, the prices were adjusting to changing user demands; yet the fairness index never fell below 0.96. No user was prevented from entering the economy nor did any user exit early. For this simulation, the pricing method was able to price link bandwidth in such a manner that lead to high utilization and an optimal distribution. Users were able to purchase link bandwidth, maximizing their QoS score and yielding a high percent GoB.

## 7 Conclusions

This paper introduced a decentralized flow control method based on microeconomics. A computer network was viewed as an economy con-

**(a)** Link 0 allocation graph.



**(b)** Fairness index graph.

**Figure 3:** Allocation and fairness index graphs.

sisting of three entities; users, Network Brokers (NB) and switches. Switches own the resources sought by users, and price their resources based on local supply and demand. A user requires these resources to maximize their individual QoS. Representing the user in the economy, the NB makes the resource purchasing decisions based on current needs of the user and prices. Users and switches act independently, which yields a decentralized flow control method. This competitive market structure encourages high utilization and Pareto optimal resource distribution. This paper also discussed how this economy properly handles network dynamics such as, users entering/exiting and VBR sources. Simulation results demonstrate the ability of the economy to successfully price link bandwidth of a network with a large number of users, each transmitting one of fifteen actual MPEG-compressed video traces. Utilization for this network was over 90% and the distribution of link bandwidth was considered optimal over 92% of the time. The price method has also been shown to perform better than standard flow control schemes [5]. This paper provided a preliminary outline and some promising experimental results. Some open questions include: wealth distribution (an issue for any economy), proper admission control

for VBR sources (especially live or interactive sources where a priori information is limited), and the possible advantage of collectively pricing a group of switches (sub-network or backbone).

# References

[1] ATM Forum Technical Committee. Traffic Management Specification. Available through `ftp://ftp.atmforum.com/pub/approved-specs/af-tm-0056.000.ps`, 1996.

[2] D. Bertsekas and R. Gallager. *Data Networks*. Prentice Hall, second edition, 1992.

[3] R. Cocchi, S. Shenker, D. Estrin, and L. Zhang. Pricing in Computer Networks: Motivation, Formulation, and Example. *IEEE/ACM Transactions on Networking*, 1(6):614 – 627, Dec 1993.

[4] D. F. Ferguson, C. Nikolaou, J. Sairamesh, and Y. Yemini. Economic Models for Allocating Resources in Computer Systems. In S. Clearwater, editor, *Market Based Control of Distributed Systems*. World Scientific Press, 1996.

[5] E. W. Fulp, M. Ott, D. Reininger, and D. S. Reeves. Congestion Pricing Flow Control for Computer Networks. Technical report, Center for Advanced Computing and Communication, April 1998.

[6] E. W. Fulp and D. S. Reeves. Dynamic Bandwidth Allocation Techniques. Technical Report TR-97/08, Center for Advanced Computing and Communication, Aug. 1997.

[7] M. Gerla and L. Kleinrock. Flow Control: A Comparative Survey. *IEEE Transactions on Communications*, 28(4):553 – 574, April 1980.

[8] R. Jain. Congestion Control and Traffic Management in ATM Networks: Recent Advances and A Survey. *Computer Networks and ISDN Systems*, Feb. 1995.

[9] H. Ji, J. Y. Hui, and E. Karasan. GoS-Based Pricing and Resource Allocation for Multimedia Broadband Networks. In *Proceedings of the IEEE INFOCOM*, pages 1020 – 1027, 1996.

[10] H. Jiang and S. Jordan. A Pricing Model for High Speed Networks with Guaranteed Quality of Service. In *Proceedings of the IEEE INFOCOM*, pages 888 – 895, 1996.

[11] J. F. Kurose and R. Simha. A Microeconomic Approach to Optimal Resource Allocation in Distributed Computer Systems. *IEEE Transactions on Computers*, 38(5):705 – 717, May 1989.

[12] S. Low and P. Varaiya. An Algorithm for Optimal Service Provisioning using Resource Pricing. In *Proceedings of the IEEE INFOCOM*, pages 368 – 373, 1994.

[13] J. K. MacKie-Mason and H. R. Varian. Pricing Congestible Network Resources. *IEEE Journal on Selected Areas in Communications*, 13(7):1141 – 1149, Sept 1995.

[14] J. Murphy and L. Murphy. Bandwidth Allocation by Pricing in ATM Networks. In *ITC*, June 1995.

[15] D. Reininger and R. Izmailov. Soft Quality-of-Service for VBR+ Video. In *Proceedings of the International Workshop on Audio-Visual Services over Packet Networks, AVSPN'97*, Sept. 1997.

[16] O. Rose. Statistical Properties of MPEG Video Traffic and Their Impact on Traffic Modeling in ATM Systems. Technical Report 101, University of Würzburg Institute of Computer Science, Feb. 1995.

[17] J. Sairamesh, D. F. Ferguson, and Y. Yemini. An Approach to Pricing, Optimal Allocation and Quality of Service Provisioning in High-speed Packet Networks. In *Proceedings of the IEEE INFOCOM*, pages 1111 – 1119, 1995.

[18] A. Takayama. *Mathematical Economics*. Cambridge University Press, 1985.

[19] L. Walras. *Elements of Pure Economics*. Richard D. Irwin, 1954. trans. W. Jaffé.